# AN ONLINE REINFORCEMENT LEARNING CONTROLLER DESIGN FOR MARS ASCENT VEHICLE

## Han Woong Bae[*]

This paper presents a neural network (NN) approximator-based online reinforcement learning (ORL) controller design for Mars Ascent Vehicle (MAV) under parametric variation and significant external disturbances. The ORL controller, which does not require any offline training, involves two NNs where an action NN produces optimal short-term control performance while a critic NN evaluates the performance of the action NN using an approximated cost function. The simulation example with comparisons against baseline Proportional-Integral-Derivative (PID) and gain scheduled pole-placement PID (GS-PP-PID) controllers show the proposed controller's effectiveness and robustness under parametric variation and high external disturbances.

## INTRODUCTION

The Mars Ascent Vehicle (MAV) concept is a part of a Mars Sample Return Campaign, where the primary objective is to return the collected samples from Martian rocks to Earth.[1] MAV is an ascent vehicle meant to deliver and eject the Orbiting Sample (OS) payload into the desired orbit.

MAV is a small and light lift-off vehicle where its target Gross Liftoff Mass (GLOM) is 400 kg. It consists of 1st and 2nd stage Solid Rocket Motors (SRM) where the first stage involves Thrust Vector Control (TVC) maneuvering while the thrust is fixed in the second stage. The vehicle also includes a hydrazine Reaction Control System (RCS) that is used to control the roll of the vehicle, whereas TVC is used to maneuver the pitch and yaw.

There are two main challenges to the design of the TVC control system: abruptly changing parameters, and unknown atmospheric conditions. The variation of the thrust of the first stage's SRM is very large where the thruster ramps up from initial thrust to peak for the initial 23 seconds and quickly drops down to less than one-half the peak. Also, the center of mass and moments of inertia change over time due to depleting solid propellant. This leads to the challenge to the design of the control system where the adaptability of the controller is needed for a robust and optimal control under the changing parameters. Furthermore, the Martian environment is not very well understood in comparison to the Earth, and access to the atmospheric survey is severely restricted. Thus, there is a high chance of risk for a significant unexpected disturbance, which could lead to the loss of stability and eventually loss of the vehicle and mission.

A simple PID control with constant gains does not sufficiently solve these challenges. The current design consideration is a gain scheduled pole placement PID (GS-PP-PID) control sys-

---

[*] Aerospace Engineer, NASA Marshall Space Flight Center, EV41/ Control System Design & Analysis.

tem, which sufficiently meets orbit requirements assuming that any unexpected disturbance does not occur. The GS-PP-PID control system uses estimated vehicle parameters and calculates the optimal PID gains using a pole placement method. This leads to significant improvements from a PID control system with constant gains; however, it does not address the outlying off-nominal conditions.

In a literature survey by Izzo et.al, the authors are confident that the machine learning techniques including deep learning and reinforcement learning will increase the level of automation as well as the performance of the Guidance Navigation and Control system.[2] This trend is continuing and growing where the AI/ML techniques are expected to be game-changers. In the remote environment where environment survey and communication are severely limited, the automation of the space systems is highly desired. Recently, an online reinforcement learning (ORL) in actor-critic structure has emerged and has shown promising performance.[3,4,5]

This ORL includes two neural networks (NN): an action NN and a critic NN. The action NN applies short-term performance measures based on output error terms, while the critic network evaluates the performance of the action based on the long-term cost function. This ORL does not require offline training, while it can quickly learn the dynamics and adapt to the changes in parameters.

Contributions of this research are the following: 1) The ORL controller for MAV does not require prior training. 2) The ORL control system does not require the knowledge of parameters and is robust against parametric change and high external disturbance. 3) The critic NN effectively eliminates steady-state error and prevents action weight gains from drifting due to the coupling effect.

## APPROACH

In this paper, three control methods are compared, PID control with constant gains, GS-PP-PID control, and ORL control. The simulation is based on MANTIS that includes sensor noises and bias, navigation latencies, and guidance steering. In this paper, pitch and yaw control using TVC during the first stage is considered.

### Rigid Body Dynamics

Assuming that the flexible body and slosh dynamics are infinitesimal, the 3-DOF equation of motion of MAV is expressed as

$$I\dot{\vec{\omega}} = -\vec{\omega} \times I\vec{\omega} + B_T\vec{u}_T \tag{1}$$

where

$$B_T = \begin{bmatrix} 0 & 0 & 0 \\ 0 & LT & 0 \\ 0 & 0 & LT \end{bmatrix} \tag{2}$$

$$\vec{u}_T = \vec{\theta}_{TVCout} \tag{3}$$

Assuming the small-angle approximation applied on the TVC gimbal angles where

$$\vec{\theta}_{TVCout} \approx \sin\left(\vec{\theta}_{TVCout}\right)$$
$$\tag{4}$$

and

$$\vec{\ddot{\theta}}_{TVCout} = \frac{1}{\omega_{TVC}^2}\vec{\dddot{\theta}}_{TVCin} + \frac{2\zeta_{TVC}}{\omega_{TVC}}\vec{\ddot{\theta}}_{TVCin} + \vec{\theta}_{TVCin} \qquad (5)$$

The variable vector $\vec{\theta}_{TVCin}$ is a gimbal angular command sent into TVC from the outer loop vehicle attitude controller and $\vec{\theta}_{TVCout}$ is an actual TVC gimbal angle as a result of an internal TVC control system controlling the gimbal angle. This internal TVC control system is assumed to have second-order dynamic with 10 Hz bandwidth and damping coefficient of 0.707 where

$$\vec{\ddot{\theta}}_{TVCout} = \frac{1}{\omega_{TVC}^2}\vec{\dddot{\theta}}_{TVCin} + \frac{2\zeta_{TVC}}{\omega_{TVC}}\vec{\ddot{\theta}}_{TVCin} + \vec{\theta}_{TVCin} \qquad (6)$$

And the vehicle attitude controller is defined by $\vec{u}_C$ and

$$\vec{\theta}_{TVCin} = \vec{u}_C \qquad (7)$$

In this paper, only the GN&C control system is considered where TVC internal control system is assumed to have second-order dynamic with 10 Hz bandwidth and damping coefficient of 0.707.

**Gain Scheduled Pole Placement PID Control**

A general form of a proportional-integral-derivative control is given by

$$\vec{u}_C = K_P\vec{e} + K_D\vec{\dot{e}} + K_I \int \vec{e}\, dt \qquad (8)$$

Where $K_P$, $K_D$ and $K_I$ are proportional, derivative and integral gains, respectively. Considering the case where gains are constant and not scheduled, the optimized values are 0.5, 0.2, and 0.2 for $K_P$, $K_D$ and $K_I$ gains, respectively.

However, the constant gains are not ideal for this application due to abruptly changing parameters such as quickly peaking magnitude of thrust, leading to degraded stability margins and control performance in extreme cases. Therefore, gain scheduled pole placement PID (GS-PP-PID) was designed to provide more robustness against the changing parameters. The parameters that are pertinent to the performance of the control system are moments of inertia, thrust moment arm, and thrust magnitude. Moments of inertia and thrust moment arm are estimated through a time-based lookup table, and the thrust magnitude is estimated through accelerometer readings. Assuming the radial symmetricity of the vehicle where $I_{yy} = I_{zz}$, the equation for a closed-loop system is set and then matched with the desired characteristic equation where the PID gains yield:

$$K_D = \frac{I_{yy}}{LT}\left(2\zeta\omega_n\lambda + \frac{1}{T_s}\right) \qquad (9)$$

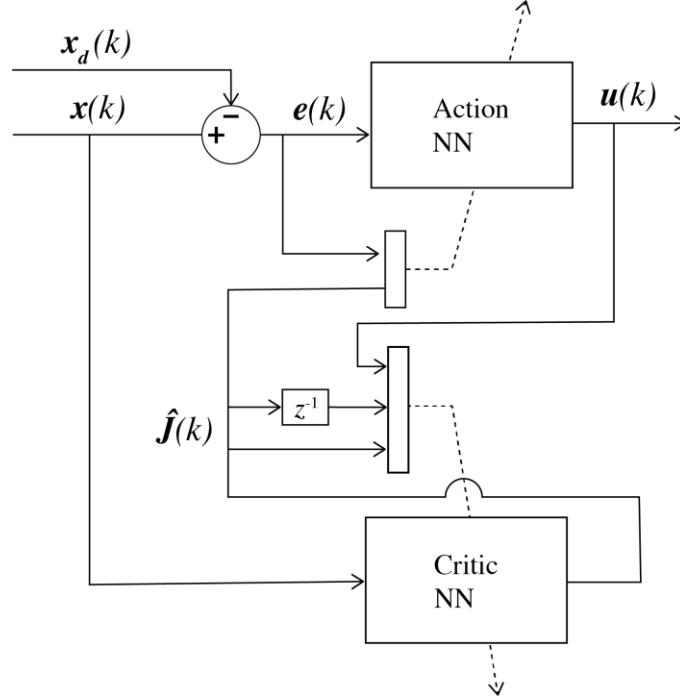$$K_P = \frac{I_{yy}}{LT}\left(\omega_n^2 + \frac{2\zeta\omega_n}{T_s}\right) \qquad (10)$$

$$K_I = \frac{I_{yy}\omega_n^2}{LTT_s} \qquad (11)$$

Where the lateral moments of inertia $I_{yy}$ is a lateral moment of inertia, $L$ is a thrust moment arm, $T$ is a magnitude of thrust, $\zeta$ is a damping coefficient, $\omega_n$ is a natural frequency, and $T_s$ is settling time where

$$T_s = \frac{4}{\zeta\omega_n} \tag{12}$$

The damping coefficient $\zeta = 1$ and natural frequency $\omega_n = 1\,\text{Hz}$ are selected for minimal overshoot and allow minimal settling time along with sufficient dynamic separation from the TVC dynamic. The derivation of the gains for GS-PP-PID is shown in Appendix A.

**Online Reinforcement Learning Control**



**Figure 1. Schematic of online reinforcement learning controller for MAV**

The online reinforcement learning control proposed for MAV's TVC control is based on Yang's method where the convergence of errors has been proven through the Lyapunov inequality and simulation results.[3] Figure 1 shows the overall schematic of the online reinforcement learning controller. The following form is a nonlinear affine system where

$$x_1(k+1) = x_2(k)$$
$$\vdots$$
$$x_n(k+1) = f\big(x(k)\big) + g\big(x(k)\big)u(k) + d(k) \tag{13}$$

In this paper, the equations of motion are expressed in continuous terms, so the conversion into the discrete term is needed to show the consistency with the theoretical stability proof via the Lyapunov Direct Method shown in Yang's paper.[3] Appendix B shows the derivation of the conversion from a continuous form into a discrete form.

The control output states that define the rotational motion of the system is given by

$$\vec{x}(k) = \vec{x}(t_k) = \left[\theta, \psi, \omega_y, \omega_z, \alpha_y, \alpha_z\right]^T \tag{14}$$

Since there are six control states, the generalized ORL control input should be six elements, where each element is designated to directly affect a corresponding element of the control states. The tracking error, $\vec{e}(k)$, at instant k is defined by

$$\vec{e}(k) = \vec{x}(k) - \vec{x}_d(k) \tag{15}$$

where $\vec{x}_d(k)$ is a desired state. The input vector of the action neural network is defined by

$$\vec{s}(k) = [\vec{e}(k)^T, 1]^T \tag{16}$$

Action NN output is given by

$$\vec{u}(k) = \hat{w}_a(k)^T \vec{\phi}'_a(v_a^T \vec{s}(k)) \tag{17}$$

where $v_a$ is an action NN first layer weight matrix with initial gains that are randomly generated and the activation function is given by $\vec{\phi}'(\vec{\beta}) = \left[\vec{\phi}(\vec{\beta}), 1\right]^T$ and

$$\vec{\phi}(\vec{\beta}) = \frac{1}{1 + e^{-\vec{\beta}/p}} \tag{18}$$

where $p$ is a tuning variable. In this paper, the activation functions for action NN and critic NN are expressed as $\vec{\phi}'_a$ and $\vec{\phi}'_c$, respectively. The action NN weight update is given by

$$\hat{w}_a(k+1) = \hat{w}_a(k) + \alpha_a \vec{\phi}'_a(v_a^T \vec{s}(k))\left(\vec{e}(k+1) - l_1\vec{e}(k) + \hat{\vec{J}}(k)\right)^T \tag{19}$$

Where $l_1$ and $\alpha_a$ are tuning variables and $\hat{\vec{J}}(k)$ is an estimated cost function approximated by critic NN. The actual cost function is represented as

$$\vec{J}(k) = \sum_{i=0}^{\infty} \gamma^i \vec{r}(k+i) \tag{20}$$

Where $\gamma$ is a discount factor ranging from 0 to 1, and $\vec{r}(k)$ is a cost function to be minimized where

$$r_i(k) = Q_i|e_i(k)|e_i(k) + R_i|u_i(k)|u_i(k), i = 1, 2, \ldots, 6 \tag{21}$$

The estimated cost function approximated by critic NN is given by

$$\hat{\vec{J}}(k) = \hat{w}_c^T(k)\vec{\phi}_c(v_c^T \vec{x}(k)) \tag{22}$$

where $v_c$ is a critic NN first layer weight matrix with initial gains that are randomly generated. The weight update law for the critic NN is given by

$$\hat{w}_c^T(k+1) = \hat{w}_c^T(k) - \alpha_c \gamma \vec{\phi}_c'(\vec{x}(k))(\gamma \hat{\vec{J}}(k) - \vec{r}(k) - \hat{\vec{J}}(k-1)) \tag{23}$$

The Bellman error is a prediction error for the critic NN and is defined as

$$\vec{e}_c(k) = \gamma \hat{\vec{J}}(k) - \vec{r}(k) - \hat{\vec{J}}(k-1)) \tag{24}$$

Lastly, the action NN output with six elements needs to be transformed into the output with two elements to control two dimensions, pitch and yaw. To do this, output variables $\{\theta, \omega_y, \alpha_y\}$ are grouped to represent pitch motion and output variables $\{\psi, \omega_z, \alpha_z\}$ are grouped to represent yaw motion. The groups are gathered in the following expression:

$$\vec{u}_c(k) = \begin{bmatrix} u_\theta + u_{\omega y} + u_{\alpha y} \\ u_\gamma + u_{\omega z} + u_{\alpha z} \end{bmatrix} = T_F \vec{u}(k) \tag{25}$$

where

$$T_F = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix} \tag{26}$$

$$\vec{u}(k) = \begin{bmatrix} u_\theta \\ u_\gamma \\ u_{\omega y} \\ u_{\omega z} \\ u_{\alpha y} \\ u_{\alpha z} \end{bmatrix} \tag{27}$$

## SIMULATION RESULTS

To demonstrate the performance of the ORL controller, it is implemented on a MANTIS (MAV Analysis Tool in Simscape) simulation tool developed by the MAV GNC team in NASA's Marshall Space Flight Center (MSFC). The flight simulation includes sensor noises and bias, navigation latencies, and guidance steering. The ORL controller is compared to the PID and GS-PP-PID controllers under different cases.
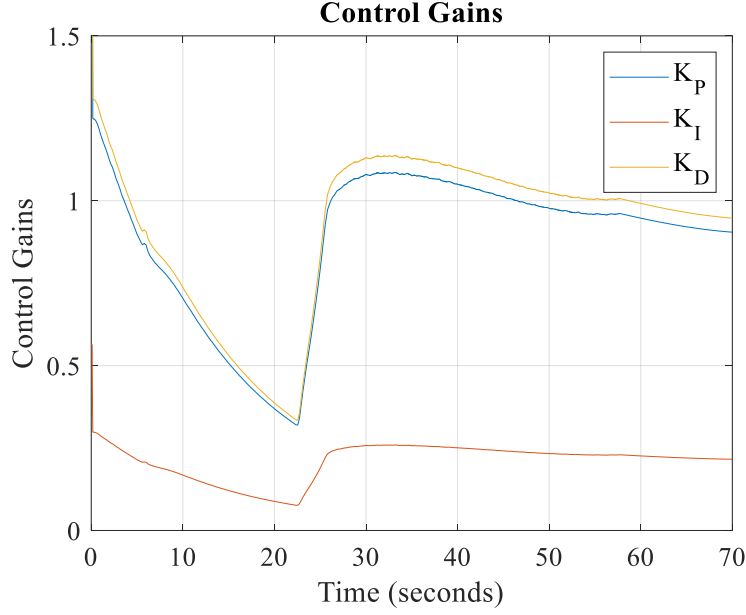
### Simulation Setup

To study the effectiveness and robustness of the proposed ORL controller, the ORL controller was tested in a simulation with different test cases. Case 1 is a nominal case both parametric uncertainties and aerodynamic external disturbance are within the expected range of dispersion. Case 2 involves an extremely high aerodynamic effect where the aerodynamic external disturbance torques are seven times the nominal. Table 1 shows the summary of cases.

**Table 1. Test Case Descriptions**

| Case # | Description |
|---|---|
| Case 1 | Nominal |
| Case 2 | High external disturbance |

Figure 2 shows the varying GS-PP-PID control gains throughout the 1st stage flight. This gain is calculated with a pole placement technique along with the estimation of parameters using a combination of look up tables and accelerometer readings. For a simple PID control with constant gains, the optimized values are 0.5, 0.2, and 0.2 for $K_P$, $K_D$ and $K_I$ gains, respectively.



**Figure 2. GS-PP-PID gains throughout the first stage TVC powered flight.**

Table 2 shows a list of control tuning variables for the ORL controller. The size of neurons in the hidden layers in action NN and critic NN are 40 and 20, respectively.
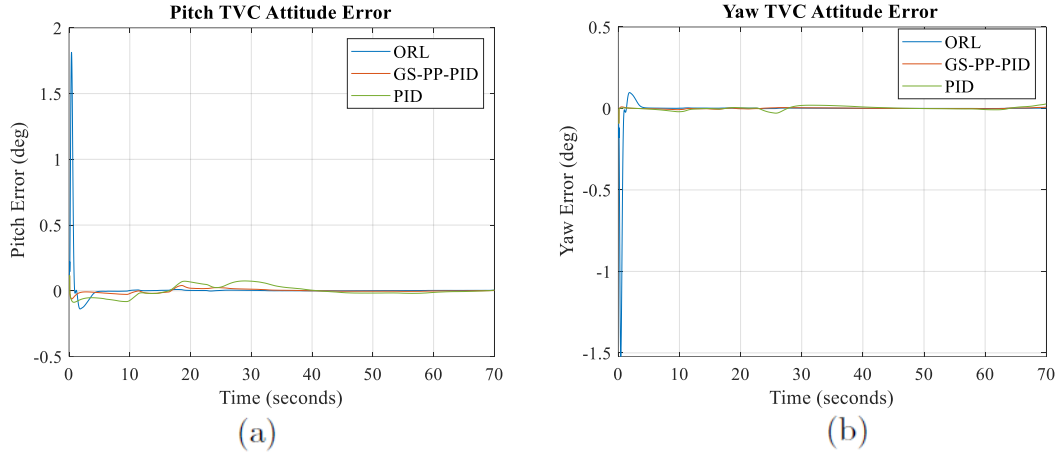
**Table 2. Control Gain Values**

| Parameter | Value |
|---|---|
| $\vec{Q}$ | $[1\ 1\ 5\ 5\ 7\ 7]$ |
| $\vec{R}$ | $[3\ 3\ 0\ 0\ 0\ 0] \times 10^5$ |
| $\alpha_a$ | $[34\ 34\ 16\ 16\ 2.1\ 2.1] \times 10^{-4}$ |
| $\alpha_c$ | $1 \times 10^{-4}$ |
| $l_1$ | 0.2 |
| $\gamma$ | 0.5 |
| $p$ | $1 \times 10^6$ |

## Results and Analysis

In this paper, time-domain attitude error results are compared. Good control performance is defined by the rate of convergence to zero as well as the ability to maintain errors close to zero throughout the flight.
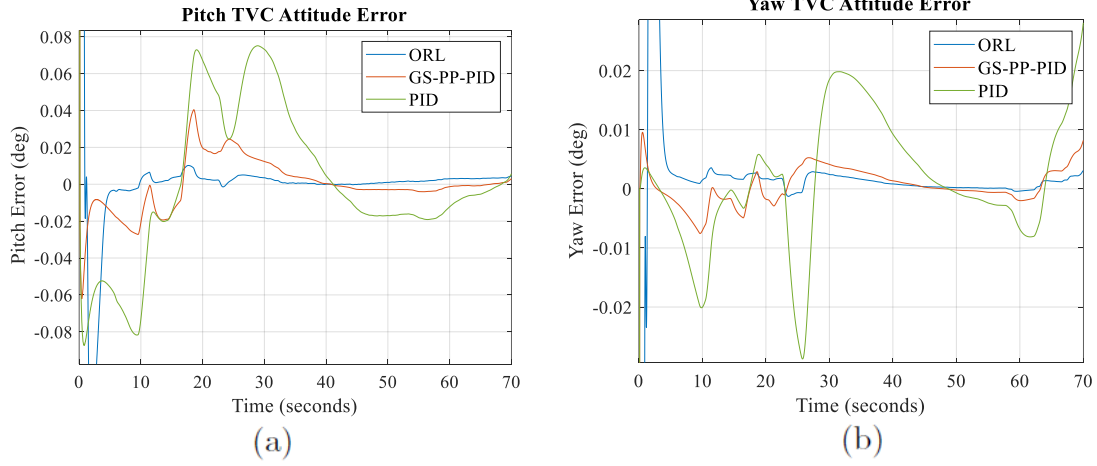
Figure 3 shows the attitude tracking errors of three controllers, ORL, GS-PP-PID, and PID, under Case 1. Note that in Figure 3, large initial transients are seen for ORL. These large transients occur due to randomization of the initial gains of the neuron weights where the controller adjusts the randomized gains to optimized gains initially, leading to high transients. For better visual comparisons, all other figures will show zoomed plots with large initial transients ignored.



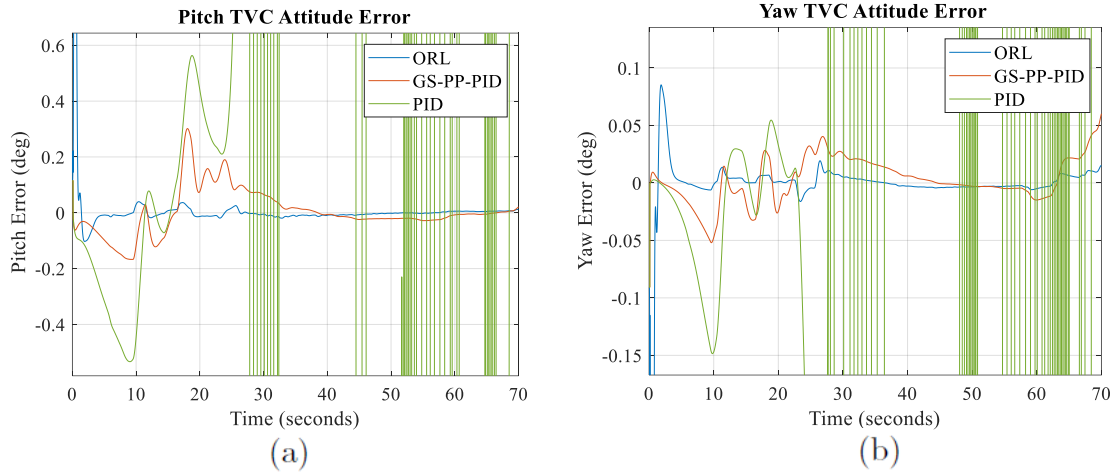**Figure 3. Vehicle (a) pitch and (b) yaw attitude errors in Case 1**

Figure 4 and 5 show the zoomed-in plots of the attitude tracking errors of three controllers in Case 1 and 2, respectively. In Figure 4, it is shown that ORL converges faster to zero as well as maintains closer to zero over time as compared to two other controllers under a nominal condition.
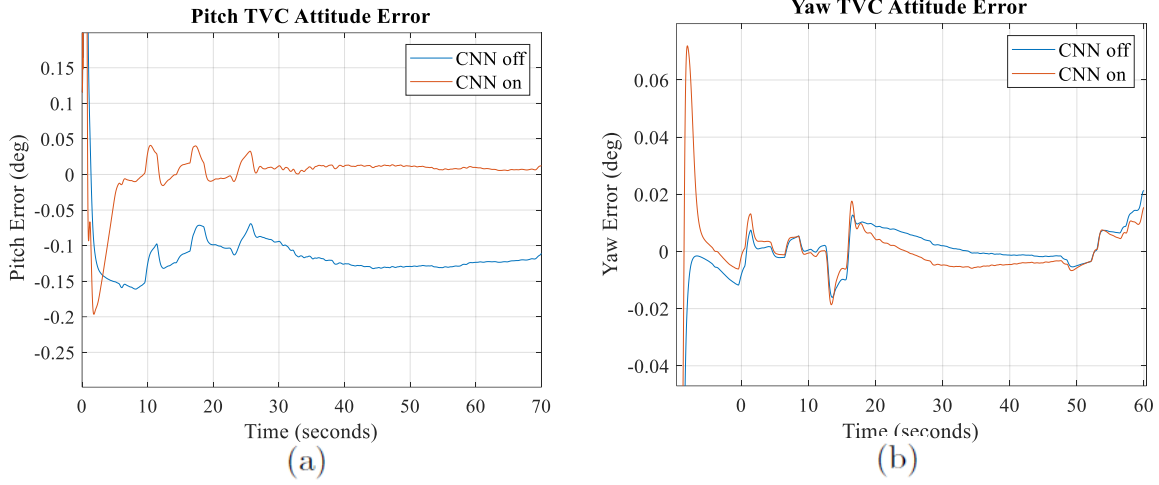
**Figure 4. Vehicle (a) pitch and (b) yaw attitude errors in Case 1 zoomed in.**

In Figure 5, as mentioned previously, Case 2 involves aerodynamic external torque disturbance that is seven times the nominal condition. This high disturbance torque has caused a loss of stability for PID control as well as marginal stability for GS-PP-PID. The ORL is able to maintain the errors at zero and shows robustness to the high disturbances.
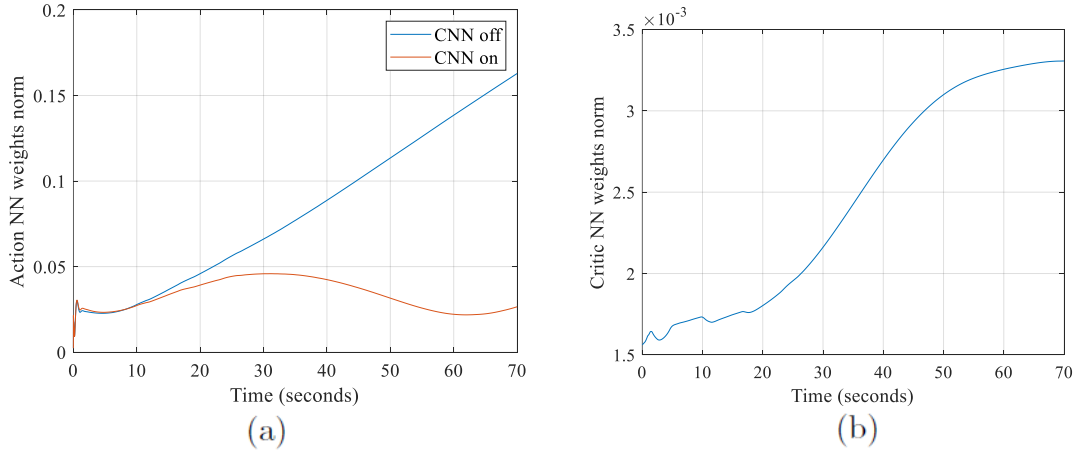


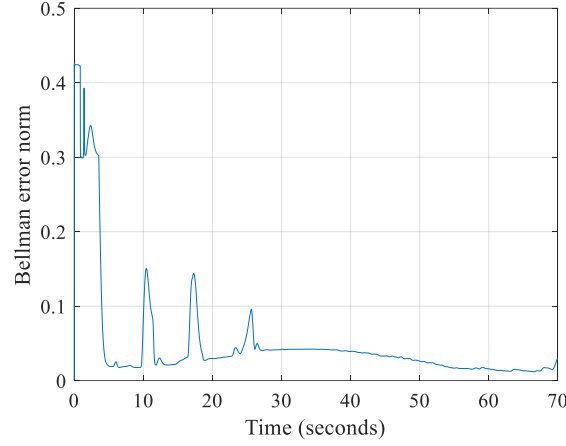**Figure 5. Vehicle (a) pitch and (b) yaw attitude errors in Case 2**

9

**Figure 6. Vehicle (a) pitch and (b) yaw attitude errors in Case 2 with Critic NN turned on and off.**

Furthermore, the effectiveness of critic NN is demonstrated in Case 2 with additional critic NN turned on and off cases. When critic NN is turned off, the pitch error is unable to converge to zero and lingers close to -0.1 deg. On the other hand, when critic NN is turned on, the pitch angle is corrected so that the error converged to zero. Also, the second advantage of critic NN is that it protects action NN weight gains from growing constantly. As seen in Figure 1(a), action NN weight increases due to coupling effects between the output states. With critic NN turned on, it penalizes action NN when output states get larger and bring the action NN weights down. Critic NN weights also stabilized toward the end of the first stage of flight.



**Figure 7. (a) Action Normed NN Weight with Critic NN turned on and off and (b) Critic Normed NN Weight**

Lastly, Figure 8 shows the normed Bellman error in Case 2. Bellman error is examined to measure the effectiveness of critic NN. As seen in the figure, Bellman error converges close to zero by the end of the flight.

10

**Figure 8. Normed Bellman error in Case 2**

## CONCLUSION

In this paper, the online reinforcement learning control has been proposed for a MAV first stage TVC control. This controller has been compared with PID control and GS-PP-PID control, where this controller has demonstrated superior robustness over two other controllers against very high aerodynamic torque disturbance. This controller is advantageous in an unpredictable Mars environment where the atmospheric survey is severely limited. With changing parameters in-flight, GS-PP-PID gains are optimized through pole placement along with real-time parameter estimation through a lookup table and accelerometer readings. The ORL, on the other hand, does not require knowledge of the parameters as it adapts to the new conditions. The future works include 1) mitigating high transients seen in the attitude errors at the beginning of flight 2) assessing the performance under multiple systematical and physical faulty conditions 3) extending ORL to guidance steering and internal TVC control for an integrated optimal solution.

## ACKNOWLEDGMENTS

## APPENDIX A: GS-PP-PID DERIVATION

Assuming geometrical symmetricity such that $I_{yy} = I_{zz}$, the control gains are assigned for the pitch and yaw are identical, the pitch channel is being considered for the derivation. The equation of motion of pitch is given by

$$I_{yy}\dot{\omega}_\theta = -\omega_\theta \times I_{yy}\omega_\theta - LF sin(\theta_{TVC}) \tag{28}$$

11

where $I_{yy}$ is a moment of inertia about the pitch axis, $\omega_\theta$ is an angular rate about the pitch axis, $L$ is a length of moment arm between the center of gravity and the thrust application point, $F$ is a thrust magnitude and $\theta_{TVC}$ is a gimbal angle in pitch.

Linearizing and simplifying the equation, the resulting equation of motion is

$$I_{yy}\ddot{\theta} = -LF\theta_{TVC} \tag{29}$$

Putting the equation of motion into the state-space form, the resulting multi-input multi-output (MIMO) equation is

$$\dot{\vec{\theta}} = A\vec{\theta} + Bu \tag{30}$$

where

$$\vec{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}, \quad \theta_1 = \int \theta\, dt, \ \ \theta_2 = \theta, \ \ \theta_3 = \dot{\theta} \tag{31}$$

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \tag{32}$$

$$B_T = \begin{bmatrix} 0 \\ 0 \\ -\dfrac{LT}{I_{yy}} \end{bmatrix} \tag{33}$$

Assuming the control input commands are zero, the control gains and law are introduced where

$$u = K\vec{\theta}, \quad K = [K_I\ K_P\ K_D] \tag{34}$$

Then the characteristic equation can be obtained and matched with the desired characteristic equation on the right side for pole placement.

$$\det\!\left(\lambda I - (A - BK)\right) = (\lambda^2 + 2\zeta\omega_n\lambda + \omega_n^2)(\lambda + \frac{1}{T_s}) \tag{35}$$

Assuming the control input commands are zero, the control gains and law are introduced where

$$K_D = \frac{I_{yy}}{LT}\left(2\zeta\omega_n\lambda + \frac{1}{T_s}\right) \tag{36}$$

$$K_P = \frac{I_{yy}}{LT}\left(\omega_n^2 + \frac{2\zeta\omega_n}{T_s}\right) \tag{37}$$

$$K_I = \frac{I_{yy}\omega_n^2}{LTT_s} \tag{38}$$

$$T_s = \frac{4}{\zeta\omega_n} \tag{39}$$

The desired control natural frequency and damping coefficient used in this study are $\omega_n = 1\ \mathrm{Hz}$ and $\zeta = 1$.

## APPENDIX B: DISCRETIZATION OF THE OUTPUT MODEL

Given control states and input states where

$$\vec{x}(k) = \vec{x}(t_k) = \left[\theta, \psi, \omega_y, \omega_z, \alpha_y, \alpha_z\right]^T \tag{40}$$

and

$$\vec{u}(t_k) = \vec{u}(k) = \begin{bmatrix} u_\theta \\ u_\gamma \\ u_{\omega y} \\ u_{\omega z} \\ u_{\alpha y} \\ u_{\alpha z} \end{bmatrix} \tag{41}$$

The continuous equation of motion is given by

$$\dot{\vec{x}}(t) = A\vec{x}(t) + B\vec{u}(t) + D \tag{42}$$

where

$$A = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \tag{43}$$

$$B = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix} \tag{44}$$

13

$$D = \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \\ d_5 \\ d_6 \end{bmatrix} \tag{45}$$

assuming that all the nonlinear part of the equation is set as disturbance and included in $D$. Then the discretized equation of motion is given by

$$\vec{x}(k+1) = \vec{f}\big(x(k)\big) + g\vec{u}(k) + \vec{d}(k) \tag{46}$$

where

$$\vec{f}\big(\vec{x}(k)\big) = (I + \Delta t A)\vec{x}(k) \tag{47}$$

$$g = \Delta t B \tag{48}$$

$$\vec{d}(k) = \Delta t D(t_k) \tag{49}$$

Where $\Delta t$ is a time-step and $t_k$ is time at instant $k$.

## REFERENCES

[1] Yaghoubi, Schnell. "Mars Ascent Vehicle Solid Propulsion Configuration", *2020 IEEE Aerospace Conference*, March 2020

[2] D. Izzo, M. Märtens, and B. Pan, A survey on artificial intelligence trends in spacecraft guidance dynamics and control. *Astrodyn.* Vol. 3, No. 4, 2019, pp. 287–299.

[3] Q. Yang and S. Jagannathan. Reinforcement learning controller design for affine nonlinear discrete-time systems using online approximators. IEEE Trans Syst Man Cybern B Cybern, 42(2):377–390, April 2012.

[4] P. He and S. Jagannathan. Reinforcement learning neural-network-based controller for nonlinear discrete-time systems with input constraints. IEEE Trans Syst Man Cybern B Cybern, 37(2):425–436, April 2007.

[5] B Xu, C Yang, and Z Shi. Reinforcement learning output feedback NN control using deterministic learning technique. IEEE Trans. Neural Netw. Learn. Syst, 25(3):635–641, March 2014.

[6] Dorf, Richard C., and Robert H. Bishop. *Modern Control Systems*. Upper Saddle River, NJ: Prentice Hall, 2010.